# High-Rate Regenerating Codes Through Layering

Birenjith Sasidharan and P. Vijay Kumar

## Abstract

In this paper, we provide explicit constructions for a class of exact-repair regenerating codes that possess a layered structure. These regenerating codes correspond to interior points on the storage-repair-bandwidth tradeoff, and compare very well in comparison to scheme that employs space-sharing between MSR and MBR codes. For the parameter set $(n, k, d = k)$ with $n < 2k - 1$, we construct a class of codes with an auxiliary parameter $w$, referred to as canonical codes. With $w$ in the range $n - k < w < k$, these codes operate in the region between the MSR point and the MBR point, and perform significantly better than the space-sharing line. They only require a field size greater than $w + n - k$. For the case of $(n, n - 1, n - 1)$, canonical codes can also be shown to achieve an interior point on the line-segment joining the MSR point and the next point of slope-discontinuity on the storage-repair-bandwidth tradeoff. Thus we establish the existence of exact-repair codes on a point other than the MSR and the MBR point on the storage-repair-bandwidth tradeoff. We also construct layered regenerating codes for general parameter set $(n, k < d, k)$, which we refer to as non-canonical codes. These codes also perform significantly better than the space-sharing line, though they require a significantly higher field size. All the codes constructed in this paper are high-rate, can repair multiple node-failures and do not require any computation at the helper nodes. We also construct optimal codes with locality in which the local codes are layered regenerating codes.

## I. Introduction

### A. Regenerating Codes

In a distributed storage system, information pertaining to a single file is distributed across multiple nodes. In the present context, a file is a collection of $K$ symbols drawn from a finite field $\mathbb{F}_q$ of size $q$. Thus a file can be represented as a $(1 \times K)$ vector over $\mathbb{F}_q$. A data collector should be able to retrieve the entire file by downloading data from any arbitrary set of $k$ nodes. Since the nodes are prone to failure, the system should be able to repair a failed node by downloading data from the remaining active nodes. In the framework of regenerating codes introduced in [1], a codeword is a $\mathbb{F}_q$-matrix of size $(\alpha \times n)$, where each column corresponds to the data stored by a single node. A failed node is regenerated by downloading $\beta \le \alpha$ symbols from any arbitrary set of $d$ nodes. These $d$ nodes are referred to as helper nodes. Since the entire file can be recovered from any arbitrary set of $k$ nodes, we must have

$$k \le d \le n - 1.$$

The total bandwidth consumed for the repair of a single node equals $d\beta$ and is termed the repair bandwidth. Thus a regenerating code is parameterized by the ordered set, $(\mathbb{F}_q, (n, k, d), (\alpha, \beta), K)$.

In the framework of regenerating codes, it is not required that the replacement node contain exactly the same symbols as did the failed node. It is only required that following regeneration, the network possess the same properties with regard to data collection and node repair as it did prior to node failure. Thus one distinguishes between functional and exact repair in a regeneration code, [1], [2], [3]. The present paper is concerned only with exact repair.

A regenerating code is said to be linear if the encoded block of $(\alpha \times n)$ matrix is a linear transformation of the $(1 \times B)$-size file vector. Linear codes offer the advantage that data recovery and node regeneration can be accomplished through low-complexity, linear operations over the field $\mathbb{F}_q$. The regenerating codes constructed in the present paper have the additional feature that no computations are needed at a helper node, a simple transfer of the contents of the helper node suffice. We will term such regenerating codes as *help-by-transfer* regenerating codes. This is distinct from the class of help-by-transfer regenerating codes discussed in [4] which have the additional feature that no computations are needed even at the replacement node.

## B. The Classical Storage-Repair-Bandwidth Tradeoff

A major result in the field of regenerating codes is the proof in [1] that uses the cut-set bound of network coding to establish that the parameters of a regenerating code must necessarily satisfy the inequality

$$K \leq \sum_{i=0}^{k-1} \min(\alpha, (d-i)\beta). \tag{1}$$

Optimal regenerating codes are those for which equality holds in (1). It turns out that for a given value of $K, k, d$, there are multiple pairs $(\alpha, \beta)$ for which equality holds in (1). It is desirable to minimize both $\alpha$ as well as $\beta$ since minimizing $\alpha$ reduces storage requirements, while minimizing $\beta$ results in a storage solution that minimizes repair bandwidth. It is not possible to minimize both $\alpha$ and $\beta$ simultaneously and thus there is a tradeoff between choices of the parameters $\alpha$ and $\beta$. The two extreme points in this tradeoff are termed the minimum storage regeneration (MSR) and minimum bandwidth regeneration (MBR) points respectively. The parameters $\alpha$ and $\beta$ for the MSR point on the tradeoff can be obtained by first minimizing $\alpha$ and then minimizing $\beta$ to obtain

$$K = k\alpha \tag{2}$$
$$\alpha = (d-k+1)\beta \tag{3}$$

Reversing the order leads to the MBR point which thus corresponds to

$$K = dk - \binom{k}{2} \tag{4}$$
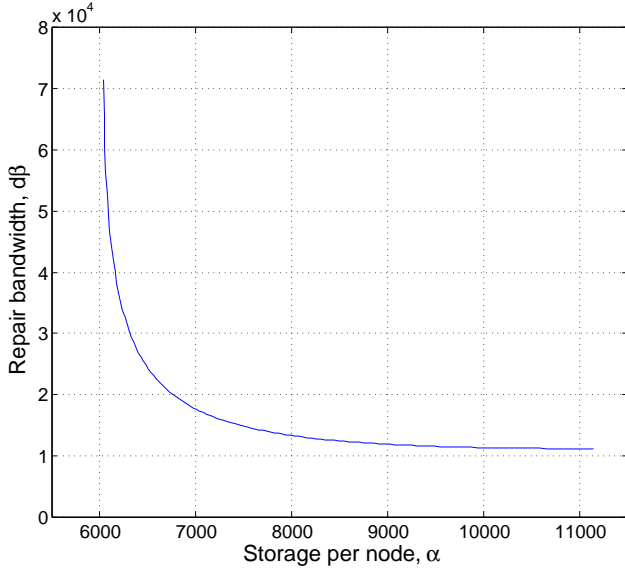$$\alpha = d\beta. \tag{5}$$



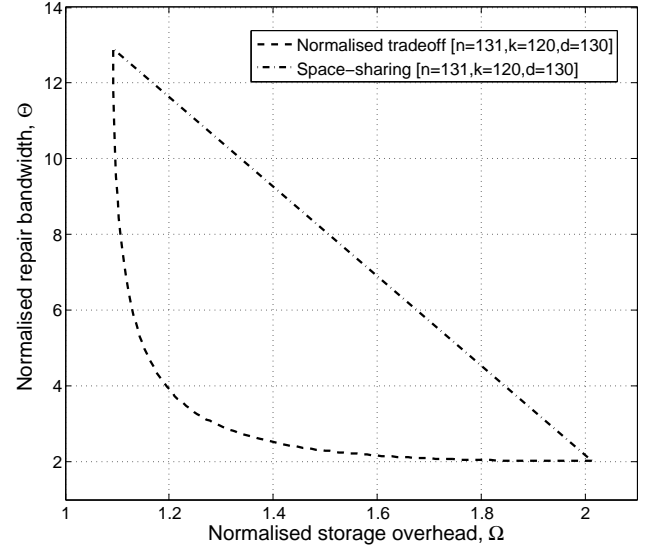Fig. 1. Storage-repair-bandwidth tradeoff. Here $[n = 131, k = 120, d = 129, B = 725360]$.



Fig. 2. Normalised storage-repair-bandwidth tradeoff. $[n = 131, k = 120, d = 129]$

The remaining points on the tradeoff will be referred to as interior points. As the tradeoff is a piecewise linear relation, there are $k$ points of slope discontinuity, corresponding to

$$\alpha = (d-p)\beta, \quad p \in \{0, 1, \cdots k-1\}.$$

Setting $p = (k-1)$ and $0$ respectively yields the MSR and MBR points respectively. Thus the remaining values of $p \in \{1, \cdots k-2\}$ correspond to interior points. The tradeoff between $\alpha$ and $d\beta$ is plotted in Fig. 1 for $(n = 131, k = 120, d = 130)$ and filesize $K = 725360$.

In [4], the authors proved that the interior points of the storage-repair-bandwidth-tradeoff cannot be achieved, under exact repair. This raises an open question as to how close one can come to the tradeoff at an interior point.

*C. Normalised Storage-Repair-Bandwidth Tradeoff*

In this subsection, we draw upon [5] and [6] to introduce a normalized version of the classical storage-repair-bandwidth tradeoff which we motivate as follows. Consider a situation where a user desires to store a file of size $K$ across $n$ nodes for a time period $T$ with each node storing $\alpha$ symbols. We follow [5] and assume a Poisson-process model of node failures under which the number of failures in time $T$ is proportional to the product of $T$ and the number of nodes $n$. We also assume that there is cost associated with both node storage as well as with repair bandwidth. The cost of storage is assumed to be proportional to the amount of data stored, i.e., to $n\alpha$. The cost of a single node-repair is taken as the amount of data download to repair a node, i.e., $d\beta$. For simplicity, we only consider the case of single-node repairs in performance comparison, although a similar analysis can be carried out for the case of multiple node failures. With this, it follows that if $\gamma(K,T)$ denotes the cost incurred to store a file of size $K$ for a time period $T$ using a particular coding scheme, then

$$\gamma(K,T) \;\; = \;\; (\gamma_B nd\beta + \gamma_S n\alpha)\, T \tag{6}$$

for some proportionality constants $\gamma_B, \gamma_S$. Hence the average cost incurred in storing one symbol for one unit of time is given by

$$\frac{\gamma(K,T)}{KT} \;\; = \;\; \gamma_B \frac{nd\beta}{K} + \gamma_S \frac{n\alpha}{K}. \tag{7}$$

We will refer to the quantities $\Omega := \frac{n\alpha}{K}$ and $\Theta := \frac{nd\beta}{K}$, as the storage overhead and normalized repair bandwidth of the code respectively. Thus the average cost is a linear combination of the normalized repair bandwidth $\Theta$ as well as the storage overhead $\Omega$. The rate $R$ of a code is the inverse of $\Omega$, i.e.,

$$R \;\; = \;\; \frac{1}{\Omega}.$$

When we set $d = n - \gamma$, and $\alpha = (d - p)\beta$, $p \in \{0, 1, \cdots, k - 1\}$, the tradeoff in (1) translates to

$$\Omega \;\; \geq \;\; \left( \frac{k}{n} - \frac{(k-p)(k-p-1)}{2n(n-\gamma)} \right)^{-1} \;\; = \;\; \Omega^* \tag{8}$$

$$\Theta \;\; \geq \;\; \frac{n-\gamma}{n-\gamma-p} \left( \frac{k}{n} - \frac{(k-p)(k-p-1)}{2n(n-\gamma)} \right)^{-1} \;\; = \;\; \Theta^* \tag{9}$$

where $\Omega^*$ and $\Theta^*$ represent the minimum possible values of $\Omega$ and $\Theta$ respectively.

A plot of of $\Theta^*$ as a function of $\Omega^*$ when $\frac{p}{n}$ is varied, is referred to as the normalised tradeoff. Unlike in the classical tradeoff, points in the modified tradeoff do not correspond to a fixed file size $K$, neither to a fixed parameter set $[n, k, d]$. The normalized tradeoff is parameterised for $\frac{k}{n}$ and $\frac{d}{n}$ where the tradeoff corresponds to the varying paramter $\frac{p}{n}$, which takes rational values bounded within $0$ and $\frac{k}{n}$. However, the plot shown in Fig. 2 is for a fixed parameter set $[n = 131, k = 120, d = 129]$.

As in [6], an asymptotic analysis of normalised storage-repair-bandwidth tradeoff can be done as $n$ scales to infinity. In this approach, the following quantities

$$\kappa \;\; = \;\; \lim_{n\to\infty} \frac{k}{n}$$

$$\theta \;\; = \;\; \lim_{n\to\infty} \frac{p}{k}, \;\; \kappa \neq 0$$

$$\Delta \;\; = \;\; \lim_{n\to\infty} \frac{d}{n}$$

are fixed as $n$ scales. If one assumes that $d = n - \gamma$, where $\gamma$ does not scale with $n$, we obtain $\Delta = 1$. Note that $\kappa \in [0\ 1]$ and, $\theta \in [0\ 1]$. Here, $\theta = 0$ correspond to the MBR point, $\theta = 1$ correspond to the MSR point, and $\theta \in (0\ 1)$ correspond to the interior points of the storage-repair-bandwidth tradeoff. In this setting, we obtain an

asymptotic version of the normalised storage-repair-bandwidth tradeoff as given below.

$$\Omega_a \geq \left( \kappa - \frac{(1-\theta)^2}{2(1-\theta\kappa)} \kappa^2 \right)^{-1} = \Omega_a^* \tag{10}$$

$$\Theta_a \geq \frac{1}{1-\theta\kappa} \left( \kappa - \frac{(1-\theta)^2}{2(1-\theta\kappa)} \kappa^2 \right)^{-1} \tag{11}$$

$$= \Theta_a^* = \frac{1}{1-\theta\kappa} \Omega_a^* \tag{12}$$

In [6], the authors have used the asymptotic analysis to study the variation of the tradeoff with respect to $\kappa$, for various points of operation $\theta$. These plots, drawn in Fig. 3, showcase the importance of regenerating codes for the interior points of the tradeoff. From Fig. 3, it follows that for any fixed storage overhead, repair bandwidth can be minimized by operating with the lowest value of $\theta$ that supports the given storage overhead. For example, if it is sufficient to build a distributed storage system with storage overhead $> 2$, then it is better to operate with $\theta = 0$, i.e., MBR point. Similarly, operating at $\theta = 1$, i.e. MSR point, is desirable only when required storage overhead is very close to 1. When the permissible storage overhead falls in the range $1 < \Omega_a < 2$, it is desirable to use codes that operate in the range $0 < \theta < 1$, i.e. in the interior region of the tradeoff.
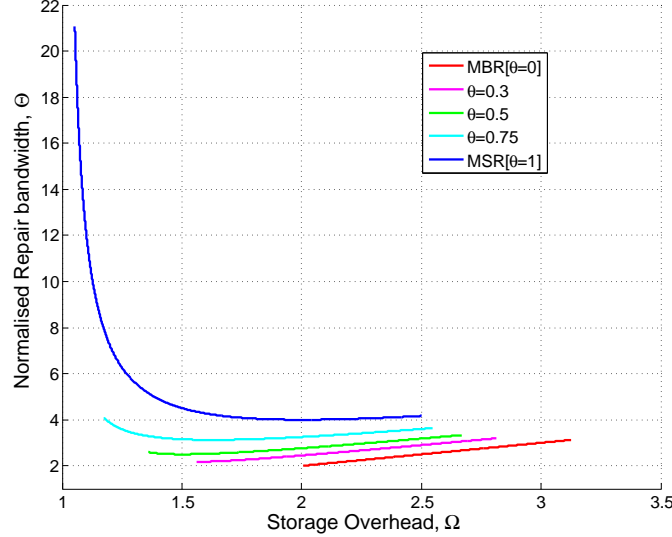


Fig. 3. Asymptotic normalised storage-repair-bandwidth tradeoff, as a function of $\kappa$, for various $\theta$.

### D. Existing Coding Schemes with Exact Repair

Several coding schemes have been proposed in the literature in the exact-repair setting. In [14], a framework to construct exact-repair optimal regenerating codes at the MBR and MSR points is provided. The framework permits the construction of MBR codes for all values for $[n, k, d]$, and of MSR codes for $d \leq 2k-3$. In [15], high-rate MSR codes with parameters $[n, k = n - 2, d = n - 1]$ are constructed using Hadamard designs. In [16], high-rate MSR codes are constructed for $d = n - 1$; here efficient node-repair is guaranteed only in the case of systematic nodes. A construction for MSR codes with $d = n - 1 \geq 2k - 1$ is presented in [17] and [18]. The construction of MSR codes for arbitrary values of $[n, k, d]$ remains an open problem, although it has been proven in [19] that exact-repair MSR codes exist for any parameter set $[n, k, d]$ as the filesize grows to infinity. In [4], a construction for a family of repair-by-transfer MBR codes is presented. The construction of regenerating codes for a functional-repair setting may be found in [20] and [18]. The nonexistence of exact-repair codes that achieve the classical storage-repair bandwidth tradeoff is proven in [4].

*E. Vector Codes*

Regenerating codes can also be viewed as vector codes. An $[n, K, d_{\min}, \alpha]$ linear *vector code* $\mathcal{C}$ over a field $\mathbb{F}_q$ is a subset of $(\mathbb{F}_q^\alpha)^n$ for some $\alpha > 1$, such that given $\mathbf{c}, \mathbf{c}' \in \mathcal{C}$ and $a, b \in \mathbb{F}_q$, $a\mathbf{c} + b\mathbf{c}'$ also belongs to $\mathcal{C}$. A codeword of a vector code is a matrix in $\mathbb{F}_q^{\alpha \times n}$, and a code symbol of a codeword is a vector in $\mathbb{F}_q^\alpha$. As a vector space over $\mathbb{F}_q$, $\mathcal{C}$ has dimension $K$, termed the scalar dimension of the code. The Hamming distance between two codewords is the number of codesymbol vectors at which they differ. In this sense, the code has minimum distance $d_{\min}$.

Associated with the vector code $\mathcal{C}$ is an $\mathbb{F}_q$-linear scalar code $\mathcal{C}^{(s)}$ of length $N = n\alpha$, where $\mathcal{C}^{(s)}$ is the collection of $(1 \times n\alpha)$ vectors obtained by vectorising each codeword matrix in some prescribed order. Given a generator matrix $G$ for the scalar code $\mathcal{C}^{(s)}$, the first code symbol in the vector code is naturally associated with the first $\alpha$ columns of $G$ and so on. We will refer to the collection of $\alpha$ columns of $G$ associated with the $i^{\text{th}}$ code symbol $\mathbf{c}_i$ as the $i^{\text{th}}$ thick column. We will refer to the columns of $G$ themselves as thin columns in order to avoid confusion, and thus there are $\alpha$ thin columns per thick column of the generator matrix.

*F. Locality*

Codes for distributed storage have been studied from other perspectives, different from the setting of regenerating codes. One prominent direction is related to codes with locality [7]. In this class of codes, a failed node is repaired by downloading entire data from a few set of nodes. Thus the property of locality allows to minimise the number of node accesses during repair. If locality-property holds only for systematic nodes, then it is referred to as information locality, and if it holds for all nodes, it is referred to as all-symbol locality. Scalar codes(i.e., $\alpha = 1$) with locality was introduced in [7], for the case of single symbol erasure, and subsequently extended in [8], for the case of multiple erasures. An upperbound on the minimum distance of the scalar code with locality was derived in above papers. Scalar codes with information locality that are optimal with respect to the aforesaid bound were constructed in [9]. Optimal scalar all-symbol local codes were constructed in [8] and [10]. Another class of codes[11], named as homomorphic self-repairing codes, constructed using linearized polynomials also turns out to be optimal scalar all-symbol codes. Recently, the concept of locality was studied for vector codes (i.e., $\alpha > 1$) in [5], [12] and [13], and thereby making this class of codes to be a comparable alternative to regenerating codes. Codes combining benefits of regenerating codes and codes with locality were constructed in [5] and [12]. In [12], the authors consider codes with all-symbol locality where the local codes are regenerating codes. Bounds on minimum distance are provided and a construction for optimal codes with MSR all-symbol locality based on linearized polynomials (rank-distance codes) is presented.

*G. Gabidulin Codes*

Let $\mathcal{G} = \{g(x) = \sum_{i=0}^{D-1} g_i x^{q^i} \mid g_i \in \mathbb{F}_q\}$ denote the set of all linearized polynomials of $q$-degree $\leq (D-1)$ over $\mathbb{F}_{q^N}$, and let $\{P_i\}_{i=1}^K$, $N \geq K \geq D$, be a collection of linearly independent elements over $\mathbb{F}_q$ in $\mathbb{F}_{q^N}$. Consider for each $g \in \mathcal{G}$, the vector $(g(P_1), g(P_2), \cdots, g(P_K))$. By representing each element $g(P_i)$ as an $N$-element vector over $\mathbb{F}_q$, we obtain an $(N \times K)$ matrix over $\mathbb{F}_q$. The resultant collection of $q^D$ matrices turns out to form a maximal rank distance (MRD) code known as the Gabidulin code [21]. In the current paper, we will in several places deal with vectors of the form $(g(P_1), \cdots, g(P_K))$, and it follows that these may also be regarded as codewords drawn from the Gabidulin code.

*H. Results*

In this paper, we first construct an $(n, k, d = k)$-regenerating code having a layered structure which we term as the canonical code $\mathcal{C}_{\text{can}}$. This code has two auxiliary parameters $w$ and $\gamma$ satisfying $w \geq 2, \gamma \geq 1, w + \gamma \leq n$ and only requires field size $q > w + \gamma$. We show how starting from a canonical code, it is possible to build a second class of layered regenerating codes with $k < d$ by making suitable use of linearized polynomial evaluations (or equivalently, codewords in the Gabidulin code) as is done in [12]. These codes will be referred to as non-canonical regenerating codes. The extension to the case $k < d$ requires however, an expansion in field size from $q$ to $q^K$ where $K$ is the scalar dimension of the underlying canonical code. These codes allows help-by-transfer repair("uncoded" repair) and are of high-rate.

We also show that the canonical code with $\gamma < w < k$ always perform better than space-sharing code. Recently, Chao et al. proposed a construction of exact-repair codes [23] using Steiner systems that achieves points better than the space-sharing line. [1] They consider constructions for $d = n - 1$. For the particular case of $d = n - 1$ and $k = n - 2$, the performance of their code is identical to the construction in the present paper when our construction is specialized to the same parameter set $d = n - 1, k = n - 2$.

Our constructions with $k = d = n - 1$ achieve an interior point of the storage-repair-bandwidth tradeoff, that is in the middle of the MSR point and the next point of slope-discontinuity. Recently, Chao [24] has characterized the optimal storage-repair-bandwidth tradeoff of $(4, 3, 3)$-exact repair codes. It turns out that the $(4, 3, 3)$-canonical code appearing in this paper also achieves the same optimal region.

Finally, we construct codes with local regeneration following the techniques in [8], [12], in which the local codes correspond to the canonical code.

*I. Performance of Codes*

The performance of this class of codes is compared against MBR and MSR codes using the normalized tradeoff. The layered codes operate in the interior region between the MSR and MBR points, and the auxiliary parameter $2 \leq w \leq k$ turns out to determine the specific interior point in the tradeoff. For a wide range of parameters $(n, k, d)$, these codes outperform codes that space-share between MSR and MBR codes [2]. Figures 4 and 5 show the respective performance of canonical codes with $(n = 61, k = 60, d = 60)$ and $(n = 61, k = 58, d = 58)$. For the case of $(n = 61, k = 60, d = 60)$, and interior point on the tradeoff between the MSR point and the next point of slope-discontinuity is achieved with $w = 59$. Achievability of interior point by canonical construction is depicted in the classical storage-repair-bandwidth plot in Fig. 7 for the parameter set $(n = 8, k = 7, d = 7)$ with auxiliary parameter $w = 6$. The performance of non-canonical layered regenerating codes with $(n = 61, k = 55, d = 60)$ is shown in Fig.6. As can be seen in plots, the codes come close to the tradeoff in terms of performance.

## II. CONSTRUCTION OF THE $(n, k, d = k)$-CANONICAL LAYERED REGENERATING CODE

In this section, we will describe the construction of a family of high-rate, $((n, k, d = k), (\alpha, \beta), K_c)$ regenerating codes indexed by two auxiliary parameters $w, \gamma$ satisfying $w \geq 2, \gamma \geq 1, w + \gamma \leq n$. The code has a layered structure, and we will have $d = k$. The code will be simply referred to as canonical code. The construction we provide in this section, assumes $(n, w + \gamma) = 1$. The general case of $(n, w + \gamma) > 1$ will be considered in the next section.

*A. Construction of the Canonical Code $\mathcal{C}$*

The construction will make use of certian other parameters derived from $w$ and $\gamma$ as defined below.

$$L = \frac{1}{n}\binom{n}{w+\gamma} \quad \text{(number of patterns)}$$

$$V = \frac{1}{w}\text{lcm}(w, w+1, \cdots, w+\gamma-1) \quad \text{(repetition factor (of each pattern))}$$

$$M = LV \quad \text{(number of layers)}$$

$$K_c = LVnw \quad \text{(scalar dimension of the canonical code)}.$$

The structure of the canonical code, can be inferred from Fig. 8 which shows the four-step process by which the incoming message vector $\underline{u}$ is encoded:

(a) The $K_c$-tuple message vector $\underline{u}$ is first partitioned into $LVn$ $w$-tuples:

$$\underline{u} \in \mathbb{F}_q^{K_c} \quad \Rightarrow \quad \left\{ \underline{u}_\tau^{(\ell,\nu)} \in \mathbb{F}_q^w \mid 1 \leq \ell \leq L, \ 1 \leq \nu \leq V, \ 0 \leq \tau \leq n-1 \right\}.$$

---

[1]Their paper appeared in the public literature only after the initial submission of our paper on arXiv.

[2]Exact-repair MSR codes are not known to exist for every value of $(n, k, d)$-tuple. Hence the achievability of the space-sharing line joining MSR point and MBR point is not always guaranteed.
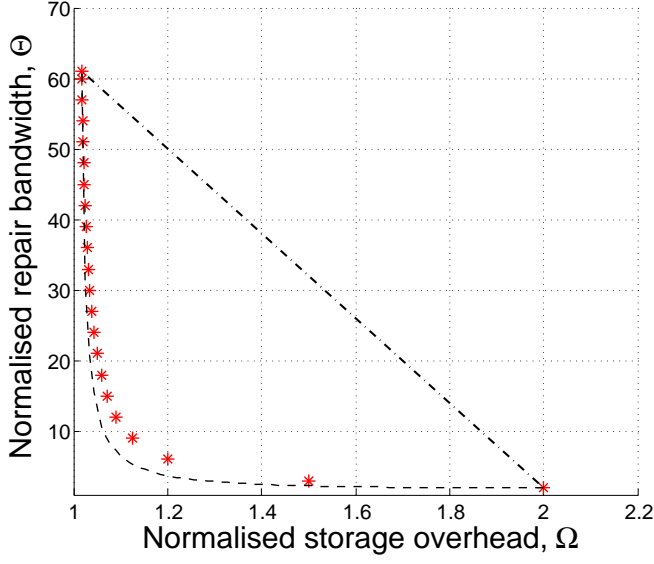
Fig. 4. Plot comparing the performance of the canonical code with $(n = 61, k = d = 60)$ for varying $w \in \{2, 5, 8, \ldots, 59\}$. The MBR point (which is the degenerate case with $w = 1$) and the MSR point are also marked. With $w = 59$, an interior point on the tradeoff between the MSR point and the next point of slope-discontinuity is achieved.
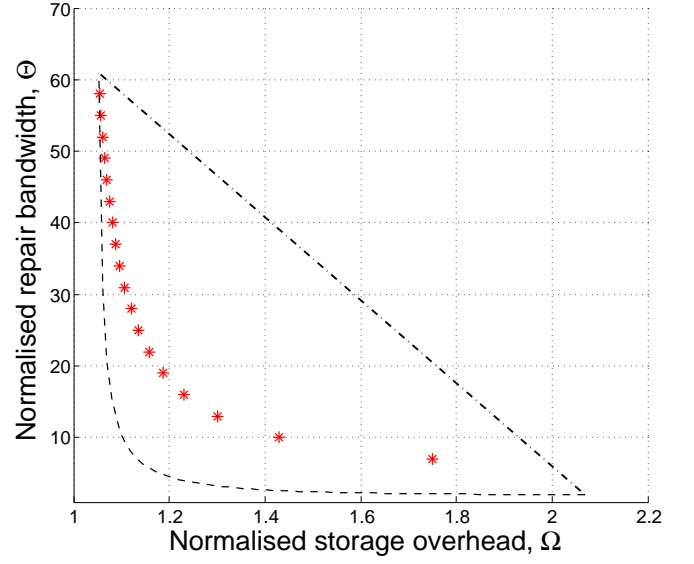


Fig. 5. Plot comparing performance of the canonical regenerating code with $(n = 61, k = 58, d = 58)$ whiel varying $w \in \{4, 7, 10, \ldots, 55\}$.
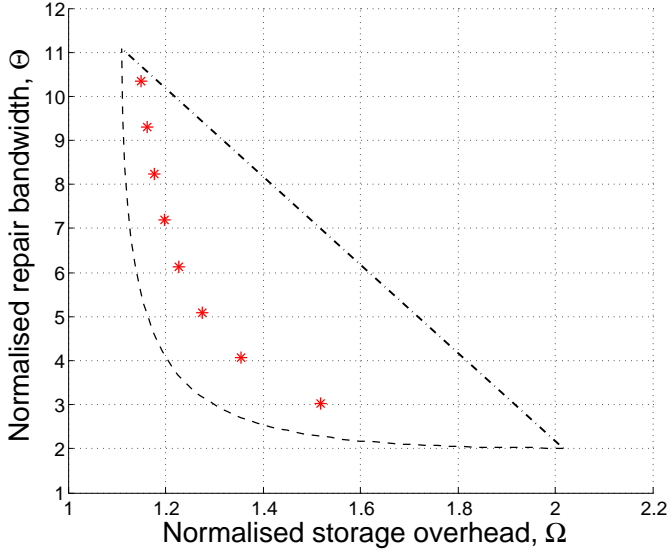


Fig. 6. Plot comparing performance of the non-canonical layered regenerating code with $(n = 61, k = 55, d = 60)$ while varying $w \in \{2, 3, \ldots, 9\}$.
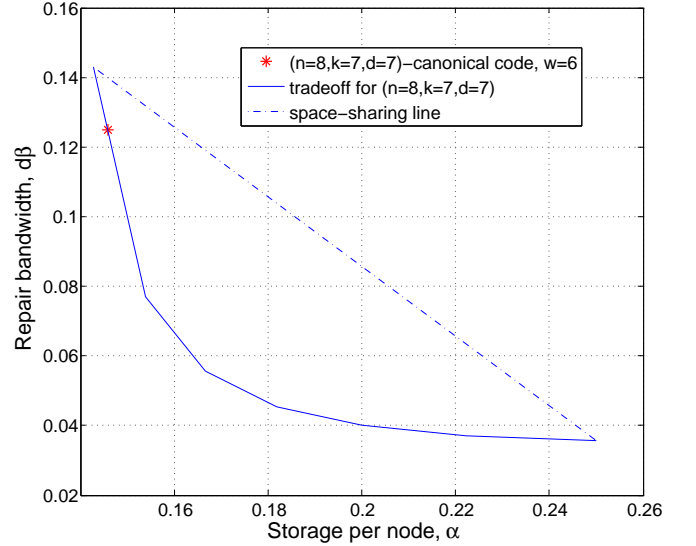


Fig. 7. $(8, 7, 7)$-canonical code with $w = 6$ achieves the interior point.

(b) Each $w$-tuple is then encoded using an $[w + \gamma, w, \gamma + 1]$ MDS code to yield $LVn$ codewords

$$\underline{u}_\tau^{(\ell,\nu)} \in \mathbb{F}_q^w \quad \Rightarrow \quad \underline{c}_\tau^{(\ell,\nu)} \in \mathbb{F}_q^{w+\gamma}.$$

(c) The collection of $n$ codewords $\{\underline{c}_\tau^{(\ell,\nu)}\}_{\tau=0}^{n-1}$ is then "threaded" to form a layer $A^{(\ell\nu)}$ of the code matrix:

$$\{\underline{c}_\tau^{(\ell,\nu)}\}_{i=0}^{n-1} \quad \Rightarrow \quad A^{(\ell,\nu)}.$$

This threading is carried out with the help of a pattern $\pi^{(\ell)}$. The nature of a pattern and the threading process are explained below in Sections II-B,II-C.

(d) The $LV$ layers are then stacked to form the code matrix

$$C \;=\; \begin{bmatrix} A^{(1,1)} \\ A^{(1,2)} \\ \vdots \\ A^{(L,V)} \end{bmatrix}.$$
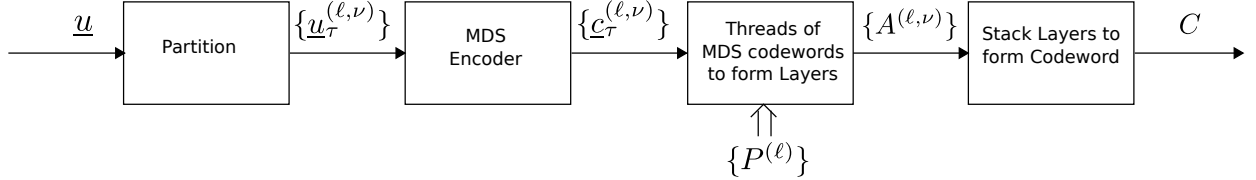


Fig. 8.   Encoder of the canonical layered regenerating code.

### B. Patterns

There are $\binom{n}{w+\gamma}$ subsets of $[n]$ that are of size $(w+\gamma)$. Let us partition these subsets into equivalence classes by declaring two elements to be equivalent if one is a cyclic shift of the other. Given our assumption that $(n, w+\gamma) = 1$, all equivalence classes will contain precisely $n$ elements and hence the number of equivalence classes is given by $L = \frac{1}{n}\binom{n}{w+\gamma}$. Let

$$\left\{ \pi^{(\ell)} \mid 1 \le \ell \le L \right\} \;=\; \left\{ (\pi_1^{(\ell)}, \pi_2^{(\ell)}, \cdots, \pi_{w+\gamma}^{(\ell)}) \mid 1 \le \ell \le L \right\},$$

be the collection of subsets obtained by selecting one subset from each equivalence class. We will assume that the elements within each of the subsets $\pi^{(\ell)}$ are ordered in ascending numerical order, i.e.,

$$\pi_1^{(\ell)} < \pi_2^{(\ell)} < \cdots < \pi_{w+\gamma}^{(\ell)}, \quad \text{for all } \ell.$$

We will associate with each such ordered subset, a collection of $n$ two-dimensional patterns, each of size $(w+\gamma) \times (w+\gamma)$. This collection includes the fundamental pattern:

$$P^{(\ell)}(0) \;=\; \left\{ \left(i, \pi_i^{(\ell)}\right) \mid 1 \le i \le w+\gamma \right\},$$

as well as its $n$ (columnar) cyclic shifts

$$P^{(\ell)}(\tau) \;=\; \left\{ \left(i, \pi_i^{(\ell)} \oplus \tau \right) \mid 1 \le i \le w+\gamma \right\}, \quad 1 \le \tau \le (n-1),$$

in which $\pi_i^{(\ell)} \oplus \tau$ is addition modulo $n$. Given a pattern $P^{(\ell)}(\tau)$ we will refer to the $(w+\gamma)$-tuple

$$\pi^{(\ell)} \oplus \tau \;=\; (\pi_1^{(\ell)} \oplus \tau,\ \pi_2^{(\ell)} \oplus \tau, \cdots,\ \pi_{w+\gamma}^{(\ell)} \oplus \tau),$$

as its (columnar) footprint. Thus the footprint of a fundamental pattern $P^{(\ell)}(0)$ is simply given by $\pi^{(\ell)}$.

### C. Threading Codewords to Form a Layer

We fix $(\ell, \nu)$ and hence describe the threading process as it applies to the $(\ell, \nu)$th layer. Consider the collection of $n$ codewords $\{\underline{c}_\tau^{(\ell,\nu)}\}_{\tau=0}^{n-1}$ associated to a layer. The symbols of the $\tau$th codeword $\underline{c}_\tau^{(\ell,\nu)}$, $0 \le \tau \le n-1$, are placed (in any arbitrary order) into the $n$ locations

$$P^{(\ell)}(\tau) \;=\; \left\{ \left(i, \pi_i^{(\ell)} \oplus \tau \right) \mid 1 \le i \le w+\gamma \right\}, \quad 1 \le \tau \le (n-1),$$

identified by the pattern $P^{(\ell)}(\tau)$. We might also refer to this codeword of this erasure code as a thread. The threading yields a $(w+\gamma \times n)$ matrix which we will denote by $A^{(\ell,\nu)}$. The threading process is illustrated in Figure 9. We then repeat this process for each layer, i.e., for all pairs $(\ell, \nu)$. Finally we vertically stack the matrices $A^{(\ell,\nu)}$ to obtain the code matrix as described above. With this the encoding process is complete.

Fig. 9. Illustrating the threading process. The top left matrix uses an $*$ to identify the elements of the two-dimensional pattern $P^{(\ell)}(0)$. The top right matrix shows the entries of a codeword $\underline{c}_0^{(\ell,\nu)}$ being inserted into the locations identified by the pattern. The bottom left shows the codeword $\underline{c}_1^{(\ell,\nu)}$ inserted into the locations identified by $P^{(\ell)}(1)$. The bottom right shows the completely filled in layer $A^{(\ell,\nu)}$.

## D. Parameters of the Canonical Code

*1) Parameters $n, \alpha$:* The parameter $n$ is simply the block length of the code $\mathcal{C}_{\text{can}}$, viewed as a vector code with symbol alphabet $\mathbb{F}_q^\alpha$. The value of $\alpha$ can be computed from nature of the construction and is given by

$$
\begin{aligned}
\alpha &= LV(w+\gamma) \\
&= \frac{1}{n}\binom{n}{w+\gamma}\frac{V_1}{w}(w+\gamma), \\
&= \frac{\text{lcm}(w, w+1, \cdots, w+\gamma-1)}{w}\binom{n-1}{w+\gamma-1}.
\end{aligned}
$$

*2) Parameters $d, \beta$:* We next note that the value of $d$ can be no less than $n-\gamma$ for otherwise, it would not be possible in some instances to repair a failed node. This follows from the fact that the symbols of each MDS code are spread across $(w+\gamma)$ distinct nodes and that to repair a failed symbol in an $[w+\gamma, w, \gamma+1]$ MDS code, one needs access to at least $w$ symbols of the codeword. Conversely, it follows that if $d = n-\gamma$, then every failed node can be repaired. We will set $d = n-\gamma$ here. It remains to establish that repair of a failed node can be accomplished by connecting to $d$ nodes and downloading a *fixed* number $\beta$ of symbols from each of the $d$ helper nodes.

It will be convenient in our analysis to assume that along with the given failed node (say node $\eta_1$), there are $\gamma-1$ other nodes (say, nodes $\eta_i$, $i = 2, 3, \ldots, \gamma$) that have also failed and that the remaining $d = n-\gamma$ nodes are acting as the helper nodes. Let us assume further, that node $h$ is one of the helper nodes. Our interest is in determining the number of symbols that need to be transferred from node $h$ to node $\eta_1$ for the purposes of node repair. We had noted earlier in describing the construction of the canonical code, that each layer $A^{(\ell,\nu)}$ of the canonical code is composed of $n$ MDS codewords $\{\underline{c}_\tau^{(\ell,\nu)}\}_{\tau=0}^{n-1}$. The codeword $\underline{c}_\tau^{(\ell,\nu)}$ is placed in the locations associated to the pattern $P^{(\ell)}(\tau)$. We will refer to the $n$ MDS codes as threads in the description below.

Node $h$ can transfer one symbol to the replacement for node $\eta_1$ iff there is a thread in some layer to which both nodes $\eta_1$ and $h$ contribute code symbols. We now break up our count according to the total number $p$ of nodes that have now failed, but which previously contributed a symbol to the erasure code thread. More specifically, we are counting the number of threads such that

- both nodes $\eta_1$ and $h$ contribute a single code symbol to that thread
- $(p-1)$ of the nodes $\{\eta_i \mid 2 \le i \le \gamma\}$ each contribute one code symbol to the thread, the remaining failed nodes do not contribute any code symbol to the thread

The total number of such threads, across all the $L$ distinct layers in the code matrix is given by

$$
\binom{\gamma-1}{p-1}\binom{n-\gamma-1}{w+\gamma-p-1}.
$$

Within the erasure code, the situation is that $p$ symbols have been erased and thus a total of $w + \gamma - p$ symbols can serve as helper nodes for node $\eta_1$ of which node $h$ is one. Since any $w$ nodes suffice to help node $\eta_1$ recover from the erasure, it suffices if node $h$ "on average" contributes a fraction

$$\frac{w}{w + \gamma - p}$$

of code symbols. We can ensure that this average is realized by calling upon the $V$ repetitions of each layer. The number $V$ has been chosen such that for all $p$,

$$\frac{w + \gamma - p}{w} \mid V.$$

Thus we can ensure that the helper node will always pass on

$$V \frac{w}{w + \gamma - p}$$

code symbols when counted across all $V$ repetitions of the corresponding erasure code. It follows that the value of $\beta$ and $d$ are given by

$$\beta = V \sum_{p=1}^{\gamma} \binom{\gamma - 1}{p - 1} \binom{n - \gamma - 1}{w + \gamma - p - 1} \frac{w}{w + \gamma - p}, \tag{13}$$

$$d = n - \gamma. \tag{14}$$

As a check, we note that each column contains $\alpha = LV(w + \gamma)$ symbols, each of which requires the transfer of $w$ symbols to enable repair. Since there are a total of $(n - \gamma)$ helper nodes, we must have that

$$\beta(n - \gamma) = w\alpha,$$

i.e., $\beta$ must equal

$$\beta = \frac{w}{(n - \gamma)} \alpha$$

$$= \frac{w}{(n - \gamma)} \frac{1}{n} \binom{n}{w + \gamma} V(w + \gamma). \tag{15}$$

It can be verified that the values for $\beta$ obtained in (16) and (17) are the same.

*3) Determining $k, K$ and Code Rate $R$:* Arguing as above, if $k < (n - \gamma)$, we will fail to decode at least one thread. Hence $k \geq (n - \gamma)$. On the other hand, by connecting to $d = (n - \gamma)$ we can recover the entire data and hence $k = d$. The scalar dimension of the code is clearly given by $K_c = LVnw$. Not surprisingly, the rate $R$ of the code is given by $\frac{w}{w + \gamma}$.

## III. $(n, k, d = k)$-CANONICAL CODE WHEN $(n, w + \gamma) \neq 1$

We consider the general case when $(w + \gamma, n) \neq 1$ and let the integer $g$ be defined by setting

$$\frac{n}{g} = (n, w + \gamma).$$

The differences in the case of $(w + \gamma, n) \neq 1$ arise out of how patterns are identified in the canonical code.

### A. Patterns

We partition as before, the $\binom{n}{w + \gamma}$ subsets of $[n]$ of size $(w + \gamma)$ into equivalence classes by declaring two subsets to be equivalent if one is a cyclic shift of the other. This time, however, different equivalence classes will be of different size. The number of elements in an equivalence class will always be of the form $gr$ with $r$ dividing $\frac{n}{g}$. Let $E(gr)$ denote the number of equivalence classes of size $gr$ and the total number of equivalence classes by $\mathcal{E}$. The values of $E(gr)$ and of $\mathcal{E}$ are given by (proof in the appendix):

$$E(gr) = \frac{1}{gr} \sum_{s \mid r} \mu(s) \binom{\frac{gr}{s}}{\frac{(w + \gamma)gr}{ns}}$$

$$\mathcal{E} = \sum_{r : gr \mid n} E(gr),$$

where $\mu(\cdot)$ denotes the Möbius function. Let

$$\left\{\pi^{(\ell)} \mid 1 \leq \ell \leq \mathcal{E}\right\} = \left\{(\pi_1^{(\ell)}, \pi_2^{(\ell)}, \cdots, \pi_{w+\gamma}^{(\ell)}) \mid 1 \leq \ell \leq \mathcal{E}\right\},$$

be the collection of subsets obtained by selecting one subset from each equivalence class. We will assume that the elements within each of the subsets $\pi^{(\ell)}$ are ordered in ascending numerical order, i.e.,

$$\pi_1^{(\ell)} < \pi_2^{(\ell)} < \cdots < \pi_{w+\gamma}^{(\ell)}, \quad \text{for all } \ell.$$

We will associate with each such subset, a collection of $n$ two-dimensional patterns, each having the same size $(w + \gamma)$. This collection includes the fundamental pattern:

$$P^{(\ell)}(0) = \left\{\left(i, \pi_i^{(\ell)}\right) \mid 1 \leq i \leq w + \gamma\right\},$$

as well as its $n$ (columnar) cyclic shifts

$$P^{(\ell)}(\tau) = \left\{\left(i, \pi_i^{(\ell)} \oplus \tau\right) \mid 1 \leq i \leq w + \gamma\right\}, \quad 1 \leq \tau \leq (n-1).$$

Given a pattern $P^{(\ell)}(\tau)$ we will refer to the $(w + \gamma)$-tuple

$$\pi^{(\ell)} \oplus \tau = (\pi_1^{(\ell)} \oplus \tau, \ \pi_2^{(\ell)} \oplus \tau, \cdots, \ \pi_{w+\gamma}^{(\ell)} \oplus \tau),$$

as its (columnar) footprint. Thus the footprint of a fundamental pattern $P^{(\ell)}(0)$ is simply given by $\pi^{(\ell)}$.

### B. Layers of the Canonical Code

Let us define

$$\begin{aligned}
L &= \sum_{r:gr|n} rE(gr) \\
V_1 &= \operatorname{lcm}(w, w+1, \cdots, w+\gamma-1) \\
V &= \frac{V_1}{w} \\
M &= LV.
\end{aligned}$$

Let us define the function $\{\omega_\ell \mid 1 \leq \ell \leq \mathcal{E}\}$ by

$$\omega_\ell = r \text{ if the pattern } \pi^{(\ell)} \text{ has period } gr.$$

It follows that there are $E(gr)$ patterns corresponding to the value $\omega_\ell = r$. Thus we can alternately express $L$ in the form

$$L = \sum_{r:gr|n} rE(gr) = \sum_{\ell=1}^{\mathcal{E}} \omega_\ell.$$

Each code matrix $C$ is composed of $LV$ vertical stacked layers, each layer corresponding to a matrix $\{A^{(\ell,\omega,\nu)} \mid 1 \leq \ell \leq \mathcal{E}, \ 1 \leq \omega \leq \omega, \ 1 \leq \nu \leq V\}$ of size $((w+\gamma) \times n)$. Thus $C$ is of the form:

$$C = \begin{bmatrix} A^{(1,1,1)} \\ \vdots \\ A^{(\ell,\omega,\nu)} \\ \vdots \\ A^{(\mathcal{E},\omega_{\mathcal{E}},V)} \end{bmatrix}.$$

The entries of the $\{A^{(\ell,\omega,\nu)}\}$ are specified below.

## C. Threading Codewords to Form a Layer

The threading process is identical to the case of $(n, w + \gamma) = 1$. We fix $(\ell, \nu)$ and the threading process in the $(\ell, \nu)$th layer is as follows. Consider the collection of $n$ codewords $\{\underline{c}_\tau^{(\ell,\nu)}\}_{\tau=0}^{n-1}$ associated to a layer. The symbols of the $\tau$th codeword $\underline{c}_\tau^{(\ell,\nu)}$, $0 \leq \tau \leq n - 1$, are placed (in any arbitrary order) into the $n$ locations

$$P^{(\ell)}(\tau) = \left\{ \left( i, \pi_i^{(\ell)} \oplus \tau \right) \mid 1 \leq i \leq w + \gamma \right\}, \quad 1 \leq \tau \leq (n-1),$$

identified by the pattern $P^{(\ell)}(\tau)$. The threading yields a $(w + \gamma \times n)$ matrix which we will denote by $A^{(\ell,\nu)}$. We then repeat this process for each layer, i.e., for all pairs $(\ell, \nu)$. It follows that a given pattern $P^{(\ell)}(\tau)$ determines the ordering of code symbols in $\omega_\ell V$ layers. In loose terms, each pattern is repeated $\omega_\ell V$ times and the parameters $\{\nu, \omega\}$ may hence be viewed as repetition parameters. The repetition factor $\omega_\ell$ that is pattern dependent, will help as we shall see, ensure a larger code rate, whereas the constant repetition factor $V$ ensures a uniform download during node repair as in Section II.

Finally we vertically stack the matrices $A^{(\ell,\nu)}$ to obtain the code matrix. This completes specification of the code matrix $C$.

## D. Parameters of the Canonical Code

*1) Parameters $n, \alpha$:* Since the number of layers change, the parameter $\alpha$ is different from the case of $(n, w+\gamma) = 1$. It is given by

$$
\begin{aligned}
\alpha &= LV(w + \gamma) \\
&= \frac{(w + \gamma) \cdot \mathrm{lcm}(w, w+1, \cdots, w+\gamma-1)}{w} \sum_{r:gr|n} rE(gr) \\
&= \frac{(w + \gamma) \cdot \mathrm{lcm}(w, w+1, \cdots, w+\gamma-1)}{wg} \sum_{r:gr|n} \sum_{s|r} \mu(s) \binom{\frac{gr}{s}}{\frac{(w+\gamma)gr}{ns}}
\end{aligned}
$$

*2) Determining Parameters $d$ and $\beta$:* Following the exact set of arguments in Sec. II-D2, we can show that $d = n - \gamma$ and $\beta$ is given by

$$\beta = V(n, w + \gamma) \sum_{p=1}^{\gamma} \binom{\gamma - 1}{p - 1} \binom{n - \gamma - 1}{w + \gamma - p - 1} \frac{w}{w + \gamma - p}. \tag{16}$$

As a check, we note that each column contains $\alpha = LV(w + \gamma)$ symbols, each of which requires the transfer of $w$ symbols to enable repair. Since there are a total of $(n - \gamma)$ helper nodes, we must have that

$$\beta(n - \gamma) = w\alpha,$$

i.e., $\beta$ must equal

$$
\begin{aligned}
\beta &= \frac{w}{(n - \gamma)} \alpha \\
&= \frac{w}{(n - \gamma)} (w + \gamma) LV. \tag{17}
\end{aligned}
$$

It can be verified that the values for $\beta$ obtained in (16) and (17) are the same.

*3) Determining $k, K$ and Code Rate $R$:* Arguing as eariler, the scalar dimension of the code is clearly given by $K_c = LVnw$. Not surprisingly, the rate $R$ of the code is given by $\frac{w}{w+\gamma}$.

## IV. CONSTRUCTION OF $(n, k < d, d)$-LAYERED REGENERATING CODE

In this section, we will describe the construction of non-canonical layered regeneration code, $\mathcal{C}_{\mathrm{lrc}}$ for general parameter set $((n, k, d), (\alpha, \beta), K)$ regenerating codes, again indexed by two auxiliary parameters $w, \gamma$ satisfying $2 \leq w < k, 1 \leq \gamma \leq (n - k)$.

## A. Construction of the non-canonical code $\mathcal{C}_{lrc}$

The non-canonical regenerating code $\mathcal{C}_{\text{lrc}}$ makes use of the canonical code code $\mathcal{C}_{\text{can}}$ as shown in Fig. 10. It also makes use of linearized polynomials along the lines of their usage in [12]. Since the construction uses the canonical code, we need to consider the case of $(n, w + \gamma) = 1$ and $(n, w + \gamma) > 1$ separately. We will consider only the case of $(n, w + \gamma) = 1$, and the general case follows accordingly.

The $K$ message symbols $\{m_i\}_{i=1}^{K}$ of $\mathcal{C}_{\text{LRC}}$ are first used to construct a linearized polynomial

$$f(x) \quad = \quad \sum_{i=1}^{K} m_i x^{q^{i-1}}.$$

The linearized polynomial is then evaluated at $K_b$ elements $\{\theta_i\}_{i=1}^{K_b}$ of $\mathbb{F}_{q^N}$ which when viewed as vectors over $\mathbb{F}_q$, are linearly independent. The resulting $K_c$ evaluations $\{f(\theta_i)\}$ are than fed as input to an encoder for the canonical $\mathcal{C}$. We set

$$u_i \quad = \quad f(\theta_i), \ 1 \le i \le K_c,$$
$$\underline{u} \quad = \quad (u_1, u_2, \ldots, u_{K_c}).$$

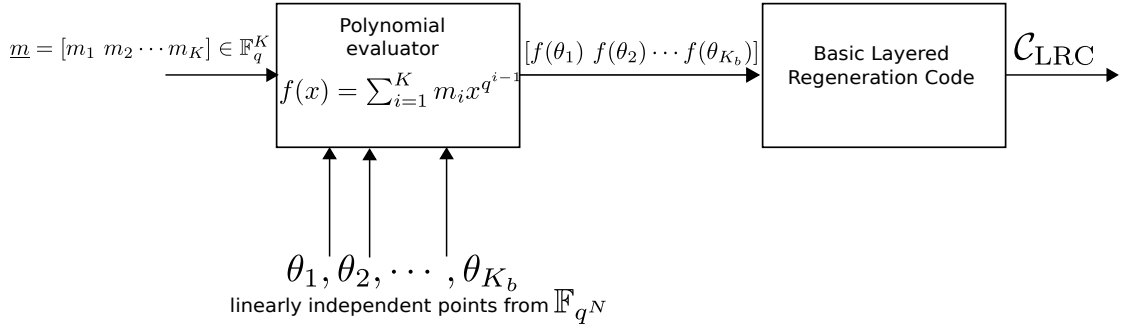The non-canonical regenerating code is the output of the canonical code to the input $\underline{u}$.



Fig. 10. Encoder of a Layered Regenerating Code.

## B. Parameters of $\mathcal{C}_{lrc}$

Clearly, the parameters $n, \alpha$ are exactly same as that of canonical code. First we proceed to relate $k$ and $K$ of $\mathcal{C}_{\text{lrc}}$. Towards that, we being with presenting a generator-matrix view point of the canonical code.

*1) Two generator matrices for the canonical code $\mathcal{C}_{can}$:* Thus far, we have described the code in terms of the structure of the codeword, viewed as a layered array. Towards determining $k$ and $K$ of the code, we now turn to a generator matrix viewpoint of the code. To obtain a generator matrix, one needs to vectorize the code matrix, thus replacing the code matrix by a vector of size $n\alpha = nLV(w + \gamma)$. The generator matrix then describes the linear relation between the $LVw$ input symbols of the canonical code $\mathcal{C}$ and the $n\alpha$ output symbols. Let us set $N_b = n\alpha$ and recall that $K_b = LVnw$. Then the generator matrix is of size $(K_b \times N_b)$.

The generator matrix is clearly dependent upon the manner in which vectorizing of the code matrix takes place. We will present two vectorization and hence, two generator matrices:

(a) From the distributed storage network point of view, the natural vectorization is one in which the $N_b$ code symbols are ordered such that the first $\alpha$ symbols correspond to the elements of the first column vector (in top-to-bottom order), of the code matrix, the second $\alpha$ symbols correspond in order, to the elements of the second column vector etc. Thus, under this vectorization, we will have that the first $\alpha$ columns of the generator matrix correspond to the first column vector of the code matrix and so on. We will refer to this as the canonical vectorization of the code. In terms of the vector-code terminology introduced earlier, each set of columns of the generator matrix corresponding to a column of the code matrix, is referred to as a thick column of the generator matrix. The code symbols associated to the $i$th thick column of the generator matrix are the code

symbols stored in the $i$th storage node. We will use $G$ to denote the generator matrix of the canonical code $\mathcal{C}$ under this vectorization.

(b) Next, consider a second vectorization of the canonical code $\mathcal{C}$ and hence, a different generator matrix. The code symbols in the code matrix of the canonical code $\mathcal{C}_{\text{can}}$ can be vectorized in such a manner that the resultant code vector is the serial concatenation of the $Mn$ codewords $\{\underline{c}_\tau^{(\ell,\nu)}\}$ of the code $\mathcal{C}_{\text{MDS}}$, each associated to a distinct message vector $\underline{u}_\tau^{(\ell,\nu)}$. Let $G_{\text{b-d}}$ denote the associated generator matrix of $\mathcal{C}_{\text{can}}$. Clearly, $G_{\text{b-d}}$ has a block-diagonal structure:

$$
G_{\text{b-d}} \;=\; \begin{bmatrix} G_{\text{MDS}} & & & \\ & G_{\text{MDS}} & & \\ & & \ddots & \\ & & & G_{\text{MDS}} \end{bmatrix}. \tag{18}
$$

Here $G_{\text{MDS}}$ denotes the generator matrix of the $[w+\gamma, w, \gamma+1]$-MDS code. It follows from this that the columns of $G_{\text{b-d}}$ associated to code symbols belonging to distinct MDS codewords are linearly independent. Also, any collection of $w$ columns of $G_{\text{b-d}}$ associated with the same $\mathcal{C}_{\text{MDS}}$ are linearly independent.

It is our intent to use the matrix $G$ for generating the canonical code $\mathcal{C}_{\text{can}}$ and the matrix $G_{\text{b-d}}$ for analysis of $\mathcal{C}_{\text{can}}$. We note that the two generator matrices $G$ and $G_{\text{b-d}}$ of the code $\mathcal{C}_{\text{can}}$ differ only in the order in which the thin columns appear.

*2) Rank Accumulation in the Matrix $G$:* The matrix $G$ has the following uniform rank-accumulation property, namely that if one selects a set $S$ containing $s$ thick columns drawn from amongst the $n$ thick columns comprising $G$, then the rank the submatrix $G|_S$ of $G$ is independent of the choice of $S$. Hence the rank of $G|_S$ may be denoted as $\rho_s$, indicating that it just depends on the value of $s$.

We now proceed to determine $\rho_s$. The value of $\rho_s$ depends on how the collection of thin columns in $S$ intersect with the blocks of $G_{\text{b-d}}$. For every thick column of $G|_S$, let us focus on a subset of thin columns corresponding to symbols from layers with a fixed value of $\nu$. We will refer to the submatrix of $G|_S$ thus obtained as $G^{(\nu)}|_S$. It is clear that rank of $G|_S$ is $V$ times the rank of $G^{(\nu)}|_S$. The intersection of $G^{(\nu)}|_S$ with blocks of $G_{\text{b-d}}$ can be sets of varying sizes, ranging from $0$ to $w+\gamma$. If the intersection is of size $p$, the rank accumulated is $\min\{p, w\}$, and thus it follows that

$$
\rho_s \;=\; V \sum_{p=1}^{\min\{s, w+\gamma\}} \binom{s}{p} \binom{n-s}{w+\gamma-p} \min\{p, w\}. \tag{19}
$$

We define the rank-accumulation profile of the matrix $G$ as the collection of integers $\{a_i\}_{i=1}^n$ given by

$$
a_1 \;=\; \rho_1 \tag{20}
$$
$$
a_i \;=\; \rho_i - \rho_{i-1}, \; 2 \le i \le n. \tag{21}
$$

It is straightforward to see that

$$
a_i \;=\; \alpha, \quad 1 \le i \le w,
$$
$$
a_i \;=\; 0, \quad k+1 \le i \le n.
$$

We will then have that

$$
\rho_s \;=\; \sum_{i=1}^{s} a_i, \quad 1 \le s \le n.
$$

*3) Parameters $k$ and $K$:* Having described the rank accumulation profile of the canonical code, we are ready to relate $k$ and $K$ of the layered code $\mathcal{C}_{\text{lrc}}$. We begin with a useful lemma.

*Lemma 4.1:* Let $k_0$ be the smallest number of thick columns of the generator matrix $G$ of the canonical code $\mathcal{C}_{\text{can}}$ such that the submatrix of $G$ obtained by selecting any $k_0$ thick columns of $G$ results in a submatrix of rank $\ge K$. Then by connecting to any $k_0$ nodes associated to the regenerating code $\mathcal{C}_{\text{lrc}}$, a data collector will be able to recover the message symbols $\{m_i\}_{i=1}^K$.

*Proof:* Let $S$ be a collection of thick of $k_0$ thick columns of the matrix $G$ such that

$$\text{Rank}(G|_S) \geq K.$$

The code symbols $(c_1, c_2, \cdots, c_n)$ of the layered regenerating code $\mathcal{C}_{\text{lrc}}$ are related to $G$ as shown below

$$(c_1, c_2, \cdots, c_n) = [f(\theta_1) \; f(\theta_2) \; \cdots f(\theta_{K_b})][G].$$

Using linearity of $f(\cdot)$, we can write this as

$$(c_1, c_2, \cdots, c_n) = f(\underbrace{[\underline{x}_1 \; \underline{x}_2 \cdots \underline{x}_{K_b}]}_{(N \times K_b)}[G]),$$

in which $\underline{x}_i \in \mathbb{F}_q^N$ is the vector representation of the element $\theta_i \in \mathbb{F}_{q^N}$. Set

$$X = [\underline{x}_1 \; \underline{x}_2 \cdots \underline{x}_{K_b}].$$

Since the $\{\underline{x}_i\}_{i=1}^{K_b}$ are linearly independent over $\mathbb{F}_q$, it follows that

$$\begin{aligned} \text{Rank}(X \cdot G|_S) &= \text{Rank}(G|_S) \\ &\geq K. \end{aligned}$$

Hence there are at least $K$ linearly independent columns in the matrix product $X \cdot G|_S$ and thus the computation $f(X \cdot G|_S)$ yields evaluations of $f(\cdot)$ in at least $K$ linearly independent points of $\mathbb{F}_{q^N}$. Since $f(\cdot)$ is of $q$-degree $(K-1)$, the coefficients of $f$ can be recovered from these $K$ evaluations. ∎

It follows from the discussion above, that in order to relate the parameters $K, k$ of $\mathcal{C}_{\text{lrc}}$, it suffices to study the canonical code $\mathcal{C}_{\text{can}}$ and determine the smallest number $k_0$ of columns of its generator matrix $G$, such that the corresponding sub matrix has rank at least $K$. But from the uniform rank accumulation property of the generator matrix $G$ of the canonical code $\mathcal{C}_{\text{can}}$ this is simply given by

$$k_0 = \min\{k \mid \rho_k \geq K\}. \tag{22}$$

Equivalently, the scalar dimension (or the filesize) of the layered regenerating code $\mathcal{C}_{\text{lrc}}$ $K$ for a given value of $k$ is given by

$$K = V \sum_{p=1}^{\min\{k, w+\gamma\}} \min\{w, p\} \binom{k}{p} \binom{n-k}{w+\gamma-p}. \tag{23}$$

*4) Parameters $d, \beta$:* From the discussion on rank accumulation profile, it follows that the scalar dimension $K$ will be strictly greater than $w\alpha$ when $w < k$, and hence we will have

$$w\alpha < K \leq k\alpha.$$

Thus, it is meaningful to have a scheme that repairs a failed node downloading $w\alpha$ symbols. Hence, we follow the same repair strategy as in the case of canonical code setting $d = n - \gamma$. We can repair any failed node downloading a fixed number $\beta$ of symbols from every helper node. The value of $\beta$ thus obtained would be

$$\beta = V \sum_{p=1}^{\gamma} \binom{\gamma-1}{p-1} \binom{n-\gamma-1}{w+\gamma-p-1} \frac{w}{w+\gamma-p}. \tag{24}$$

It can also be checked that $(n-\gamma)\beta = w\alpha$.

## C. Some Remarks on the parameters of $\mathcal{C}_{lrc}$

The following remarks on parameters of $\mathcal{C}_{lrc}$ are worth mentioning.

*Remark 1:* From the description in Sec. I-A, it is clear that every regenerating code must satisfy

$$(d - k + 1)\beta \;\; \leq \;\; \alpha \;\; \leq \;\; d\beta. \tag{25}$$

Since layered regenerating codes have

$$d\beta \;\; = \;\; w\alpha \; ,$$

we must have

$$k \;\; \geq \;\; 1 + d\left(\frac{w-1}{w}\right).$$

A lowerbound on $k$ imposes only a lower limit on the rate, and hence the above constraint does not come along with any penalty.

*Remark 2:* In the construction, we have assumed the auxiliary parameter $w$ to be greater than 1 because $w$ turns out to be the dimension of the common erasure code. Nevertheless, we can consider the extreme case of $w = 1$, where the erasure code becomes a trivial repetition code. In addition, let us set $\gamma = 1$, and hence $w + \gamma = 2$. Then for all odd valued $n$, $(n, w + \gamma) = 1$ and hence in that case, we have

$$
\begin{aligned}
L \;\; &= \;\; \frac{n-1}{2}, \\
V \;\; &= \;\; 1, \\
\alpha \;\; &= \;\; n - 1, \\
d\beta \;\; &= \;\; \alpha.
\end{aligned}
$$

The code thus obtained is structurally similar to the repair-by-transfer MBR codes and differs only in that the underlying MDS code present in the construction of the repair-by-transfer MBR codes in [4] is replaced here by an MDS code that is constructed using linearized polynomials.

*Remark 3:* If the linearized polynomial of $q$-degree $(K-1)$ used in the construction of $\mathcal{C}_{lrc}$ is replaced by an (ordinary) polynomial of degree $(K-1)$, then one can then still go onto to obtain a regenerating code. While this code will have smaller field size, it will however, have lesser rate in comparison to the code $\mathcal{C}_{lrc}$ constructed here.

## V. On the Optimality of the canonical code

In this section, we state two results pertaining to the performance of the canonical code against the storage-repair-bandwidth tradeoff. The first result shows that for any $(n, k, d = k)$ parameter set, we can construct canonical codes that performs better than what the space-sharing code achieves. In the second, we will establish the achievability of an interior point in the storage-repair-bandwidth tradeoff by an exact-repair code when $d = k = n - 1$. The interior point we achieve is on the line-segment joining the MSR point and the next point of slope-discontinuity, where the non-achievability results established in [4] does not apply. Both these results follow immediately from simple calculations.

*Lemma 5.1:* The $(n, k, d = k)$-canonical code operates at an $(\alpha, d\beta)$-point that lies between the MSR and MBR points, and performs better than the code that space-shares the MSR and MBR point, whenever $\gamma < w < k$.

*Proof:* For any regenerating code with $d = k$, we must have

$$\frac{\alpha}{d} \leq \beta \leq \alpha.$$

Since $w + \gamma \leq n$, we must have $w \leq n - \gamma = d = k$. Furthermore, $\beta = \frac{w}{d}\alpha$ for the canonical code. Thus it follows that code operate at a point between the MSR and MBR point.

From [4], we can express the space-sharing line in the form,

$$d\beta \;\; = \;\; \frac{d(2K - k\alpha)}{k(d - k + 1)}. \tag{26}$$

When $d = k$, it reduces to

$$d\beta = 2K - k\alpha.$$

For the canonical code, we have

$$K = \left(\frac{w}{w+\gamma}\right) n\alpha,$$

and hence, for it to perform better than the space-sharing code, we must have,

$$w\alpha < 2\left(\frac{w}{w+\gamma}\right) n\alpha - k\alpha.$$

It can be verified that the above condition holds whenever $(k-w)(w-\gamma) > 0$, which is true when $\gamma < w < k$. ∎

*Corollary 5.2:* When $n < 2k - 1$, there exist exact-repair $(n, k, d = k)$-regenerating codes that operate between the MSR and the MBR point performing better than the space-sharing line.

*Proof:* An integer value of $w$ satisfying $\gamma < w < k$ can be found when $n < 2k - 1$. The statement follows from that. ∎

*Lemma 5.3:* The $(n, n - 1, n - 1)$-canonical code achieves an interior point of the storage-repair-bandwidth tradeoff, that lies between the MSR point and the next point of slope-discontinuity specified by,

$$\alpha = (d - (k - 2))\beta - \left(\frac{k-2}{k-1}\right)\beta. \tag{27}$$

*Proof:* The results in [4] imply that the rank accumulation profile of a linear optimal regenerating code must satisfy,

$$a_p^* = \begin{cases} \min\{\alpha, (d - p + 1)\beta\}, & 1 \le p \le k \\ 0 & k < p \le n \end{cases}$$

$$= \begin{cases} \alpha & 1 \le p \le \left\lfloor d - \left(\frac{\alpha}{\beta}\right) + 1 \right\rfloor \\ (d - p + 1)\beta & \left\lfloor d - \left(\frac{\alpha}{\beta}\right) + 1 \right\rfloor < p \le k \\ 0 & k < p \le n \end{cases}$$

Thus a linear code is an optimal regenerating code if and only if it satisfies the above rank accumulation profile. For a regenerating code with $d\beta = w\alpha$, we calculate

$$a_p^* = \begin{cases} \alpha, & 1 \le p \le \left\lfloor d\left(\frac{w-1}{w}\right) + 1 \right\rfloor \\ w\alpha - \frac{w(p-1)\alpha}{d} & \left\lfloor d\left(\frac{w-1}{w}\right) + 1 \right\rfloor < p \le k \\ 0 & k < p \le n \end{cases} \tag{28}$$

Now consider the $(n, n - 1, n - 1)$-canonical code. This means i.e., $\gamma = 1$, and then it follows from (21) that the rank accumulation profile of the code,

$$a_p = \begin{cases} \alpha, & 1 \le p \le w \\ \alpha - \binom{p-1}{w}, & w + 1 \le p \le k \\ 0. & k < p \le n \end{cases} \tag{29}$$

For (28) and (29) to match, it is also necessary to check that

$$\left\lfloor d\left(\frac{w-1}{w}\right) + 1 \right\rfloor = w.$$
$$\Leftrightarrow \quad d \in \{w, w + 1\}$$

If we choose $w = d - 1$, we obtain

$$a_p = a_p^* = \alpha, \quad 1 \le p \le w = k - 1$$

Furthermore, we must check that $a_k = a_k^*$.

$$a_k = \alpha - \binom{k-1}{k-1} = \alpha - 1, \tag{30}$$

$$a_k^* = w\alpha - \frac{w^2\alpha}{w+1}$$
$$= \frac{w\alpha}{w+1}. \tag{31}$$

For the canonical code, since $(n, w+1) = (w+2, w+1) = 1$,

$$
\begin{aligned}
\alpha &= \frac{1}{n}\binom{n}{w+1}(w+1) \\
&= \frac{1}{w+2}\binom{w+2}{w+1}(w+1) \\
&= w+1
\end{aligned}
$$

Thus it follows that,

$$a_k = a_k^* = w,$$

showing that $(n, n-1, n-1)$-canonical code with $w = d-1$ is an optimal regenerating code. Since the accumulation profile has values $a_p = \alpha, 1 \le p \le (k-1)$ and $a_k = \alpha - 1$, it achieves a point between the MSR point and the next point of slope-discontinuity on the tradeoff. It can be calculated that the interior point thus achieved is specified by (27). ∎

## VI. CODES WITH CANONICAL-CODE-LOCALITY

In this section, we will briefly describe how it is possible to construct codes with locality in which each of the local codes is the canonical (layered regenerating) code $\mathcal{C}_{\text{can}}$. The same technique can also be used to generate codes with locality in which the local codes are the layered regeneration codes $\mathcal{C}_{\text{lrc}}$.

### A. Locality in Vector Codes

Let $\mathcal{C}$ be an $[n, K, d_{\min}, \alpha]$ vector code over a field $\mathbb{F}_q$, possessing a $(K \times n\alpha)$ generator matrix $G$. The $i^{\text{th}}$ code symbol, $\mathbf{c}_i$, is said to have (exact) $(r, \delta)$ locality, $\delta \ge 2$, if it is possible to puncture the code in coordinates corresponding to a set of indices $S$ with $i \in S$, such that the punctured code $\mathcal{C}|_S$ has length $r + \delta - 1$, and minimum distance $\delta$. The code $\mathcal{C}$ is said to have $(r, \delta)$ all-symbol locality if all code symbols have $(r, \delta)$ locality. The codes obtained through puncturing will be called local codes. Our interest here is in the construction of a code with exact, all-symbol locality, whose local codes correspond to the canonical code $\mathcal{C}_{\text{can}}$ introduced in Section II.

The property of locality allows to minimise the number of node accesses during node-repair. The concept of locality was introduced in [7] for scalar codes for single erasures. Subsequently it was generalised to multiple erasures and later to vector codes; see [8], [13], [5] and [12], [11]. Codes combining benefits of regenerating codes and codes with locality are constructed in [12], [5] and [22].

### B. Code Construction

Let $t \ge 2$ and $\{\phi_i\}_{i=1}^{tK_c}$ a collection of elements in $\mathbb{F}_{q^N}$ and let $\{\underline{\phi}_i\}$ denote the representation of the $\{\phi_i\}$ as elements of $\mathbb{F}_q^N$. Given a message vector $[m_1, m_2 \dots, m_K]^T$, we construct the linearized polynomial

$$h(x) = \sum_{i=1}^{K} m_i x^{q^{i-1}}, \quad m_i \in \mathbb{F}_{q^N}, \quad N \ge tK_c,$$

and form the $tK_c$-tuple $[h(\phi_1), h(\phi_2), \dots, h(\phi_{tK_c})]^T$. This evaluation vector is then partitioned into $t$ evaluation vectors each counting $K_c$ components which are then fed to $t$ respective encoders for the canonical code. The corresponding outputs of these encoders are then concatenated to form the desired codeword. It can be shown that the resultant code is optimal in terms of having the best possible minimum distance for the given scalar dimension.

APPENDIX

*Lemma A.1:* For $r$ such that $gr|n$, the number of equivalence classes of size $gr$ is given by

$$E(gr) = \frac{1}{gr} \sum_{s|r} \mu(s) \binom{\frac{gr}{s}}{\frac{(w+\gamma)gr}{ns}}.$$

In particular, the total number of equivalence classes $\mathcal{E}$ is given by

$$\mathcal{E} = \sum_{r|\frac{n}{g}} E(gr).$$

*Proof:* For $r$ such that $gr \mid n$, let $f_1(r)$ denote the number of equivalence classes of size less than or equal to $gr$. Then $f_1(r)$ is given by

$$f_1(r) \;\;=\;\; \binom{gr}{\frac{(w+\gamma)gr}{n}}.$$

Let $f_2(r)$ denote the number of patterns having size equal to $gr$. Then we have,

$$f_1(r) = \sum_{s|r} f_2(s)$$

and by Möbius inversion, we obtain

$$f_2(r) = \sum_{s|r} f_1\left(\frac{r}{s}\right) \mu(s), \tag{32}$$

where $\mu$ is the Möbius function. Thus the number of equivalence classes of size $gr$ is given by

$$
\begin{aligned}
E(gr) \;\;&=\;\; \frac{1}{gr} f_2(r) \\
&=\;\; \frac{1}{gr} \sum_{s|r} f_1\left(\frac{r}{s}\right) \mu(s) \\
&=\;\; \frac{1}{gr} \sum_{s|r} \mu(s) \binom{\frac{gr}{s}}{\frac{(w+\gamma)gr}{ns}}.
\end{aligned}
$$

The result for the total number of equivalence classes follows immediately. ∎

REFERENCES

[1] A. Dimakis, P. Godfrey, Y. Wu, M. Wainwright, and K. Ramchandran, "Network coding for distributed storage systems," *Information Theory, IEEE Transactions on*, vol. 56, no. 9, pp. 4539–4551, 2010.

[2] Wu, Y. and Dimakis, A.G., "Reducing repair traffic for erasure coding-based storage via interference alignment," in *Information Theory, 2009. ISIT 2009. IEEE International Symposium on*. IEEE, 2009, pp. 2276–2280.

[3] Rashmi, KV and Shah, N.B. and Kumar, P.V. and Ramchandran, K., "Explicit construction of optimal exact regenerating codes for distributed storage," in *Communication, Control, and Computing, 2009. Allerton 2009. 47th Annual Allerton Conference on*. IEEE, 2009, pp. 1243–1249.

[4] N. Shah, K. Rashmi, P. Vijay Kumar, and K. Ramchandran, "Distributed Storage Codes With Repair-by-Transfer and Nonachievability of Interior Points on the Storage-Bandwidth Tradeoff," *Information Theory, IEEE Transactions on*, vol. 58, no. 3, pp. 1837–1852, march 2012.

[5] Kamath, G.M. and Prakash, N. and Lalitha, V. and Kumar, P.V., "Codes with Local Regeneration," *arXiv preprint arXiv:1211.1932*, 2012.

[6] G. M. Kamath and P. V. Kumar, "Regenerating codes: A reformulated storage-bandwidth trade-off and a new construction," in *Communications (NCC), 2012 National Conference on*, feb. 2012, pp. 1–5.

[7] P. Gopalan, C. Huang, H. Simitci, and S. Yekhanin, "On the Locality of Codeword Symbols," *Electronic Colloquium on Computational Complexity (ECCC)*, vol. 18, p. 100, 2011.

[8] N. Prakash, G. M. Kamath, V. Lalitha, and P. V. Kumar, "Optimal linear codes with a local-error-correction property," in *Information Theory Proceedings (ISIT), 2012 IEEE International Symposium on*, july 2012, pp. 2776–2780.

[9] C. Huang, M. Chen, and J. Li, "Pyramid codes: Flexible schemes to trade space for access efficiency in reliable data storage systems," in *Network Computing and Applications, 2007. NCA 2007. Sixth IEEE International Symposium on*. IEEE, 2007, pp. 79–86.

[10] J. Han and L. A. Lastras-Montano, "Reliable Memories with Subline Accesses," in *Information Theory, 2007. ISIT 2007. IEEE International Symposium on*, june 2007, pp. 2531–2535.

[11] F. Oggier and A. Datta, "Self-repairing homomorphic codes for distributed storage systems," in *INFOCOM, 2011 Proceedings IEEE*. IEEE, 2011, pp. 1215–1223.

[12] A. S. Rawat, O. O. Koyluoglu, N. Silberstein, and S. Vishwanath, "Optimal locally repairable and secure codes for distributed storage systems," *CoRR*, vol. abs/1210.6954, 2012.

[13] D. S. Papailiopoulos and A. G. Dimakis, "Locally repairable codes," in *Information Theory Proceedings (ISIT), 2012 IEEE International Symposium on*, july 2012, pp. 2771–2775.

[14] K. Rashmi, N. Shah, and P. Kumar, "Optimal Exact-Regenerating Codes for Distributed Storage at the MSR and MBR Points via a Product-Matrix Construction," *Information Theory, IEEE Transactions on*, vol. 57, no. 8, pp. 5227–5239, aug. 2011.

[15] D. S. Papailiopoulos, A. G. Dimakis, and V. R. Cadambe, "Repair Optimal Erasure Codes through Hadamard Designs," *CoRR*, vol. abs/1106.1634, 2011.

[16] I. Tamo, Z. Wang, and J. Bruck, "Zigzag Codes: MDS Array Codes with Optimal Rebuilding," *CoRR*, vol. abs/1112.0371, 2011.

[17] C. Suh and K. Ramchandran, "Exact-repair MDS code construction using interference alignment," *Information Theory, IEEE Transactions on*, vol. 57, no. 3, pp. 1425–1442, 2011.

[18] N. Shah, K. Rashmi, P. V. Kumar, and K. Ramchandran, "Interference Alignment in Regenerating Codes for Distributed Storage: Necessity and Code Constructions," *Information Theory, IEEE Transactions on*, vol. 58, no. 4, pp. 2134–2158, april 2012.

[19] V. R. Cadambe, S. A. Jafar, H. Maleki, K. Ramchandran, and C. Suh, "Asymptotic Interference Alignment for Optimal Repair of MDS codes in Distributed Data Storage," 2012, preprint available online at http://www.mit.edu/~viveck/resources/Research/asymptotic_storage.pdf, Accessed 25 Oct 2012.

[20] Hu, Y. and Lee, P. P. C. and Shum, K. W., "Analysis and Construction of Functional Regenerating Codes with Uncoded Repair for Distributed Storage Systems," *arXiv preprint arXiv:1208.2787*, 2012.

[21] E. M. Gabidulin, "Theory of Codes with maximum rank distance," *Information Transmission, Problems of*, vol. 21, no. 7, pp. 1–12, Jul 1985.

[22] Kamath, Govinda M. and Prakash, N. and Lalitha, V. and Kumar, P. Vijay and Natalia Silberstein and Ankit Singh Rawat and Onur Ozan Koyluoglu and Sriram Vishwanath, "Explicit MBR All–Symbol Locality Codes," *arXiv preprint*, 2013.

[23] Chao Tian and Vaneet Aggarwal and Vinay V. Vaishampayan, "Exact-Repair Regenerating Codes Via Layered Erasure Correction and Block Designs," *arXiv preprint*, 2013.

[24] Chao Tian, "Rate Region of the $(4, 3, 3)$ Exact-Repair Regenerating Codes," *arXiv preprint*, 2013.